

# Team-Tamil@LT-EDI-RANLP2023: Automatic Detection of Hope Speech in Bulgarian Language using Embedding Techniques

Rahul Ponnusamy<sup>1</sup>, Malliga Subramaniam<sup>2</sup>, Sajeetha Thavareesan<sup>3</sup>, Ruba Priyadharshini<sup>4</sup>

<sup>1</sup> Insight SFI Research Centre for Data Analytics, University of Galway, Ireland, India

<sup>2</sup> Kongu Engineering College, Tamil Nadu, India <sup>3</sup> Eastern University, Sri Lanka

<sup>4</sup> Department of Mathematics, Gandhigram Rural Institute-Deemed to be University, Tamil Nadu, India

rahulponnusamy160032@gmail.com

mallinishanth72@gmail.com

sajeethas@esn.ac.lk, rubapriyadharshini.a@gmail.com

## Abstract

Many people may find motivation in their lives by spreading content on social media that is encouraging or hopeful. Creating an effective model which helps in accurately predicting the target class is a challenging task. The problem of Hope speech identification is dealt with in this work using machine learning and deep learning methods. In this paper, we present the description of the system submitted by our team (Team-Tamil) to the Hope Speech Detection for Equality, Diversity, and Inclusion (HSD-EDI) LT-EDI-RANLP 2023 shared task for the Bulgarian language. The main goal of this shared task is to identify the given text into the Hope speech or Non-Hope speech category. The proposed method using H2O deep learning model with MPNet embeddings and achieved the second rank for the Bulgarian language with the Macro F1 score of 0.69.

## 1 Introduction

One of the remarkable human characteristics is hope, which enables a person to envision future events and the diversity of outcomes that may be anticipated (Snyder, 1994). These visions have a significant effect on a person's emotions, behaviors, and mental state, despite the fact that the desired outcome has a much-reduced likelihood of occurring. Hope is essential to the well-being, recuperation, and restoration of human existence. Greater optimism is consistently associated with improved academic, athletic, and physical health, psychological adjustment, and psychotherapeutic outcomes. Hope theory is similar to learned optimism, optimism, self-efficacy, and self-esteem theories (Snyder, 2002; Chakravarthi, 2022a,b; García-Baena et al., 2023).

People are able to freely express their opinions on numerous social networks today, which has a

significant impact on human existence (B and Varsha, 2022; Subramanian et al., 2022). The significant characteristics of social media, including rapid dissemination, low cost, accessibility, and anonymity, have increased the popularity of social media platforms such as Instagram and Twitter (B and A, 2021; Chakravarthi et al., 2023a). Despite the numerous advantages of using OSNs, a growing body of evidence suggests that an increasing number of malicious actors are exploiting these networks to disseminate hate speech and cause harm to others (Chakravarthi et al., 2023b; Santhiya et al., 2023). In addition, social media platforms provide a profound comprehension of people's behaviors and are important sources for Natural Language Processing (NLP)-related scientific research (Chakravarthi et al., 2022a; Shanmugavadivel et al., 2022; Chakravarthi, 2023).

Examining people's expressed levels of hope on social media is therefore believed to be an essential factor in determining their overall happiness (Kumaresan et al., 2023; Subramanian et al., 2023). This type of research can shed light on the progression of goal-directed activities, resilience in the face of adversity, and the mechanisms underlying acclimation to both positive and negative life changes.

In this paper, we present the work carried out on HSD-EDI - LT-EDI-RANLP 2023 in Bulgarian language. The main goal of the shared task is to categorize the comments into Hope speech or Non-Hope speech class. To solve this problem, our team (Team-Tamil) presents the approach based on an embedding technique using MPNet (Song et al., 2020) sentence transformer and deep learning technique with H2O (Candel et al., 2016) deep learning model. Our approach achieved the second rank with a macro F1 score of 0.69 in the Bulgarian language.

The rest of the paper is structured as follows:

Table 1: Data statistics

| Dataset     | Hope speech | Non-Hope speech | Total |
|-------------|-------------|-----------------|-------|
| Train       | 223         | 4448            | 4671  |
| Development | 75          | 514             | 589   |
| Test        | 150         | 449             | 599   |

Section 2 provides a brief overview of the studies related to the Hope speech detection problem, Section 3 provides a full explanation of the dataset, Section 4 shows the proposed approach, and Section 5 provides the findings of the experiments that were conducted and followed by the conclusion in Section 6.

## 2 Related Work

Two shared task was released by Chakravarthi and Muralidaran (2021) and Chakravarthi et al. (2022b). Mahajan et al. (2021) used RoBERTa for Hope Speech detection in English and XLM-RoBERTa for Hope Speech detection in Tamil and Malayalam, two low-resourced Dravidian languages. Their performance in classifying text into hope-speech, non-hope, and not-language is demonstrated. Their methodology was rated first in English(F1 = 0.93), first in Tamil(F1 = 0.61), and third in Malayalam(F1 = 0.83). Junaida and Ajees (2021) used Deep learning-based context-aware string embeddings for word representations and Recurrent Neural Network(RNN) and aggregated document embeddings for text representation. The authors examined and contrasted each language’s three models using diverse methodologies. Their approach outperforms baselines, and English, Tamil, and Malayalam models beat baselines by 3%, 2%, and 11%, respectively. Pre-processing and transfer-learning models enable the trials.

Dowlagar and Mamidi (2021) used pre-trained multilingual-BERT with convolution neural networks for English, Malayalam-English, and Tamil-English code-mixed datasets, and they ranked first, third, and fourth. S et al. (2021) used transformer models, mBERT for Tamil and Malayalam and BERT for English, yielded weighted average F1-scores of 0.46, 0.81, and 0.92 for Tamil, Malayalam, and English, respectively. Vijayakumar et al. (2022) used BERT to do this work, and their model ranked first in Kannada, second in Malayalam, third in Tamil, and sixth in English for the hope speech 2022 shared task. B et al. (2022) used m-BERT, MLNet, BERT, XLMRoberta, and

XLM\_MLM to identify and classify them. BERT and m-BERT had the highest weighted F1-scores of 0.92, 0.71, 0.76, 0.87, and 0.83 for English, Tamil, Spanish, Kannada, and Malayalam, respectively.

Our study varies from the previous research in which we used MPNet and doc2vec(Le and Mikolov, 2014) for creating the embedding of the comments. For detecting the Hope speech, we employed H2O deep learning model.

## 3 Dataset and Task Description

We participated in the HSD-EDI task in LT-EDI-RANLP 2023. The main challenge of this task is to create a model that automatically detects whether the comment is a Hope speech or a Non-Hope speech. For this task, the organizers provided the dataset with annotated labels for Bulgarian, English, Hindi, and Spanish languages. Out of these four languages, we worked only on the Bulgarian language. The comments are annotated with two labels: True and False for the Hope speech and Non-Hope speech. The in-depth details of the dataset are provided in (Chakravarthi, 2020). In the first phase, a training and development set was released for creating the model. In the second phase, test sets were released only with the comments to make the prediction with the model that was created with the training set in the first phase. We need to submit the prediction that took from the test set to the organizers. In the last phase, the test set with the labels will be released test set with labels to know the performance of the model. The statistics of the datasets are shown in Table 1.

## 4 Methodology

We conducted an in-depth analysis of the Bulgarian Hope speech dataset utilizing a range of classifiers, from basic machine learning methods to powerful deep learning algorithms. The manner we carried out our research is outlined below. We utilized the scikit-learn<sup>1</sup> library for building machine learning algorithms. We used the h2o library to implement

<sup>1</sup><https://scikit-learn.org>

Table 2: List of parameters that are used for creating embedding using doc2vec

| Parameters  | Values |
|-------------|--------|
| dm          | 0      |
| vector_size | 300    |
| negative    | 5      |
| min_count   | 1      |
| alpha       | 0.065  |
| min_alpha   | 0.065  |

Table 3: The table shows the model results on the Development set using doc2vec(ACC: Accuracy, MAC\_P: Macro Precision, MAC\_R: Macro Recall, MAC\_F1: Macro F1, WEL\_P: Weighted Precision, WEL\_R: Weighted Recall, WEL\_F1: Weighted F1)

| MODELS     | ACC  | MAC_P | MAC_R | MAC_F1      | WEL_P | WEL_R | WEL_F1 |
|------------|------|-------|-------|-------------|-------|-------|--------|
| <b>H2O</b> | 0.84 | 0.63  | 0.62  | <b>0.62</b> | 0.83  | 0.84  | 0.84   |
| <b>DT</b>  | 0.83 | 0.58  | 0.56  | 0.57        | 0.81  | 0.83  | 0.82   |
| <b>LR</b>  | 0.87 | 0.72  | 0.53  | 0.53        | 0.84  | 0.87  | 0.83   |
| <b>RF</b>  | 0.87 | 0.44  | 0.50  | 0.47        | 0.76  | 0.87  | 0.81   |
| <b>SVM</b> | 0.87 | 0.44  | 0.50  | 0.47        | 0.76  | 0.87  | 0.81   |

Table 4: The table shows the model results on the Development set using MPNet(ACC: Accuracy, MAC\_P: Macro Precision, MAC\_R: Macro Recall, MAC\_F1: Macro F1, WEL\_P: Weighted Precision, WEL\_R: Weighted Recall, WEL\_F1: Weighted F1)

| MODELS     | ACC  | MAC_P | MAC_R | MAC_F1      | WEL_P | WEL_R | WEL_F1 |
|------------|------|-------|-------|-------------|-------|-------|--------|
| <b>H2O</b> | 0.83 | 0.65  | 0.68  | <b>0.66</b> | 0.85  | 0.83  | 0.84   |
| <b>DT</b>  | 0.83 | 0.56  | 0.53  | 0.54        | 0.80  | 0.83  | 0.81   |
| <b>RF</b>  | 0.87 | 0.44  | 0.50  | 0.47        | 0.76  | 0.87  | 0.81   |
| <b>SVM</b> | 0.87 | 0.44  | 0.50  | 0.47        | 0.76  | 0.87  | 0.81   |
| <b>LR</b>  | 0.87 | 0.44  | 0.50  | 0.47        | 0.76  | 0.87  | 0.81   |

the deep neural network. We used Google Colab-  
 oratory<sup>2</sup> to train the models because of its user  
 interface and quicker access to GPU resources.

For detecting the hope speech from the text,  
 Firstly, we remove noise from data in order to im-  
 prove data quality for improved performance and  
 remove URLs and Unhelpful expressions(terms  
 that start with @). Secondly, we converted the  
 cleaned text to feature vectors using doc2vec and  
 MPNet. The doc2vec<sup>3</sup> is also known as para-  
 graph vectors, an unsupervised method for learn-  
 ing fixed-length feature representations from texts  
 with varying lengths, such as paragraphs, sen-  
 tences, and documents. Each sentence is repre-  
 sented by a dense vector. We set up parameters with  
 dm=0(training algorithm with a distributed bag of  
 words),vector\_size=300(dimensionality of the fea-  
 ture vectors), negative=5 for negative sampling,

<sup>2</sup><https://colab.research.google.com>

<sup>3</sup><https://radimrehurek.com/gensim/models/doc2vec.html>

and the remaining parameters are listed in Table  
 2. MPNet is a cutting-edge pre-training technique  
 that inherits the benefits of BERT and XLNet while  
 avoiding their drawbacks. We get embedding from  
 the text using the pretrained model(‘all-mpnet-base-  
 v2’ from sentence transformer<sup>4</sup>). Thirdly, we ex-  
 perimented with the traditional machine learning  
 techniques, namely, Logistic Regression(LR), Ran-  
 dom Forest(RF), Decision Tree(DT), and Support  
 Vector Machine(SVM) Classifier using the scikit-  
 learn library and the Deep learning technique using  
 H2O framework<sup>5</sup> with rectifier as activation func-  
 tion and [100,100] of hidden layer size. The model  
 is trained for ten epochs.

In the next section, we explained the perfor-  
 mance of the models.

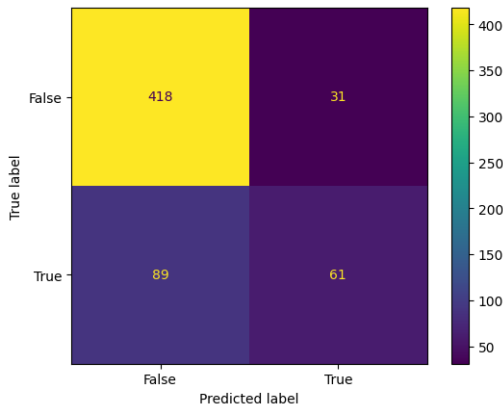
<sup>4</sup><https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

<sup>5</sup><https://github.com/h2oai/h2o-3/tree/master>

Table 5: This table shows the result on the Test set with the H2O model using MPNet embeddings

| H2O_MPNet          | Scores      |
|--------------------|-------------|
| Accuracy           | 0.80        |
| Macro Precision    | 0.74        |
| Macro Recall       | 0.67        |
| Macro F1           | <b>0.69</b> |
| Weighted Precision | 0.78        |
| Weighted Recall    | 0.80        |
| Weighted F1        | 0.78        |

Figure 1: Confusion matrix of the H2O deep learning model with MPNet on the test set. Support value for False(Non-Hope speech) is 449 and True(Hope speech) is 150.



## 5 Results and discussion

In this section, we discussed the outcomes of the experiments of the model that we used. The performance of the models is evaluated with Accuracy, Macro Precision, Macro Recall, Macro F1, Weighted Precision, Weighted Recall, and Weighted F1 score. There are 449 samples of Non-Hope speech and 150 samples of Hope speech. We used four machine learning models and one deep learning model with two embedding techniques: doc2vec and MPNet. Among all other models, H2O deep learning model performed well with both doc2vec and MPNet embeddings on the development set with macro F1 scores of 0.62 and 0.66, respectively. The results of the models with doc2vec and MPNet are shown in Table 3 and Table 4. We selected the top-performing model on the development set, that is H2O deep learning model with MPNet embeddings, to make predictions on the test set. The final leaderboard results revealed that the proposed methodology ranked in second place in the Bulgarian language with a Macro F1-score of 0.69. The results of the test set are shown

in Table 5. The confusion matrix in Figure 1 displays the right prediction as 418 out of 449 is false and 61 out of 150 is true.

## 6 Conclusion

Social media platforms have evolved into a forum for people to discuss their thoughts, successes, achievements, and errors. Social networking members leave comments on various types of content. Positive words can assist in increasing confidence and sometimes push you to be strong in difficult situations. This article describes our model that was submitted to the HSD-EDI - LT-EDI-RANLP 2023 competition. We used the H2O framework to build a deep neural network for categorization. We utilized MPNet from the sentence transformer for creating embeddings. Our proposed method comes in second place with a weighted F1 score of 0.69 which is above the baselines. To improve performance further, the model can be fine-tuned with model architecture as well as by doing hyperparameter tuning.

## References

- Bharathi B and Agnusimmaculate Silvia A. 2021. *SS-NCSE\_NLP@DravidianLangTech-EACL2021: Offensive language identification on multilingual code mixing text*. In *Proceedings of the First Workshop on Speech and Language Technologies for Dravidian Languages*, pages 313–318, Kyiv. Association for Computational Linguistics.
- Bharathi B, Dhanya Srinivasan, Josephine Varsha, Thenmozhi Durairaj, and Senthil Kumar B. 2022. *SS-NCSE\_NLP@LT-EDI-ACL2022: hope speech detection for equality, diversity and inclusion using sentence transformers*. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 218–222, Dublin, Ireland. Association for Computational Linguistics.
- Bharathi B and Josephine Varsha. 2022. *SSNCSE\_NLP@TamilNLP-ACL2022: Transformer based ap*



- proach for detection of abusive comment for Tamil language. In *Proceedings of the Second Workshop on Speech and Language Technologies for Dravidian Languages*, pages 158–164, Dublin, Ireland. Association for Computational Linguistics.
- Arno Candel, Viraj Parmar, Erin LeDell, and Anisha Arora. 2016. Deep learning with h2o. *H2O. ai Inc*, pages 1–21.
- Bharathi Raja Chakravarthi. 2020. Hopeedi: A multilingual hope speech detection dataset for equality, diversity, and inclusion. In *Proceedings of the Third Workshop on Computational Modeling of People’s Opinions, Personality, and Emotion’s in Social Media*, pages 41–53.
- Bharathi Raja Chakravarthi. 2022a. Hope speech detection in youtube comments. *Social Network Analysis and Mining*, 12(1):75.
- Bharathi Raja Chakravarthi. 2022b. Multilingual hope speech detection in english and dravidian languages. *International Journal of Data Science and Analytics*, 14(4):389–406.
- Bharathi Raja Chakravarthi. 2023. Detection of homophobia and transphobia in youtube comments. *International Journal of Data Science and Analytics*, pages 1–20.
- Bharathi Raja Chakravarthi, Adeep Hande, Rahul Ponnusamy, Prasanna Kumar Kumaresan, and Ruba Priyadharshini. 2022a. How can we detect homophobia and transphobia? experiments in a multilingual code-mixed setting for social media governance. *International Journal of Information Management Data Insights*, 2(2):100119.
- Bharathi Raja Chakravarthi, Manoj Balaji Jagadeeshan, Vasanth Palanikumar, and Ruba Priyadharshini. 2023a. Offensive language identification in dravidian languages using mpnet and cnn. *International Journal of Information Management Data Insights*, 3(1):100151.
- Bharathi Raja Chakravarthi and Vigneshwaran Muralidaran. 2021. Findings of the shared task on hope speech detection for equality, diversity, and inclusion. In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 61–72, Kyiv. Association for Computational Linguistics.
- Bharathi Raja Chakravarthi, Vigneshwaran Muralidaran, Ruba Priyadharshini, Subalalitha Cn, John McCrae, Miguel Ángel García, Salud María Jiménez-Zafra, Rafael Valencia-García, Prasanna Kumaresan, Rahul Ponnusamy, Daniel García-Baena, and José García-Díaz. 2022b. Overview of the shared task on hope speech detection for equality, diversity, and inclusion. In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 378–388, Dublin, Ireland. Association for Computational Linguistics.
- Bharathi Raja Chakravarthi, Ruba Priyadharshini, Shubanker Banerjee, Manoj Balaji Jagadeeshan, Prasanna Kumar Kumaresan, Rahul Ponnusamy, Sean Benhur, and John Philip McCrae. 2023b. Detecting abusive comments at a fine-grained level in a low-resource language. *Natural Language Processing Journal*, 3:100006.
- Suman Dowlagar and Radhika Mamidi. 2021. EDIOne@LT-EDI-EACL2021: Pre-trained transformers with convolutional neural networks for hope speech detection. In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 86–91, Kyiv. Association for Computational Linguistics.
- Daniel García-Baena, Miguel Ángel García-Cumbreras, Salud María Jiménez-Zafra, José Antonio García-Díaz, and Rafael Valencia-García. 2023. Hope speech detection in spanish: The lgbt case. *Language Resources and Evaluation*, pages 1–28.
- MK Junaida and AP Ajees. 2021. Ku\_nlp@ lt-edi-eacl2021: a multilingual hope speech detection for equality, diversity, and inclusion using context aware embeddings. In *Proceedings of the first workshop on language technology for equality, diversity and inclusion*, pages 79–85.
- Prasanna Kumar Kumaresan, Rahul Ponnusamy, Elizabeth Sherly, Sangeetha Sivanesan, and Bharathi Raja Chakravarthi. 2023. Transformer based hope speech comment classification in code-mixed text. In *Speech and Language Technologies for Low-Resource Languages*, pages 120–137, Cham. Springer International Publishing.
- Quoc Le and Tomas Mikolov. 2014. Distributed representations of sentences and documents. In *International conference on machine learning*, pages 1188–1196. PMLR.
- Khyati Mahajan, Erfan Al-Hossami, and Samira Shaikh. 2021. TeamUNCC@LT-EDI-EACL2021: Hope speech detection using transfer learning with transformers. In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 136–142, Kyiv. Association for Computational Linguistics.
- Arunima S, Akshay Ramakrishnan, Avantika Balaji, Thenmozhi D., and Senthil Kumar B. 2021. ssn\_diBERTsity@LT-EDI-EACL2021:hope speech detection on multilingual YouTube comments via transformer based approach. In *Proceedings of the First Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 92–97, Kyiv. Association for Computational Linguistics.
- S. Santhiya, P. Jayadharshini, and S. V. Kogilavani. 2023. Transfer learning based Youtube toxic comments identification. In *Speech and Language Technologies for Low-Resource Languages*, pages 220–230, Cham. Springer International Publishing.

- Kogilavani Shanmugavadivel, Sai Haritha Sampath, Pramod Nandhakumar, Prasath Mahalingam, Malliga Subramanian, Prasanna Kumar Kumaresan, and Ruba Priyadharshini. 2022. [An analysis of machine learning models for sentiment analysis of Tamil code-mixed data](#). *Computer Speech Language*, 76:101407.
- C Richard Snyder. 2002. Hope theory: Rainbows in the mind. *Psychological inquiry*, 13(4):249–275.
- Charles Richard Snyder. 1994. *The psychology of hope: You can get there from here*. Simon and Schuster.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. Mpnet: Masked and permuted pre-training for language understanding. *Advances in Neural Information Processing Systems*, 33:16857–16867.
- Malliga Subramanian, Ramya Chinnasamy, Prasanna Kumar Kumaresan, Vasanth Palanikumar, Madhoora Mohan, and Kogilavani Shanmugavadivel. 2023. Development of multi-lingual models for detecting hope speech texts from social media comments. In *Speech and Language Technologies for Low-Resource Languages*, pages 209–219, Cham. Springer International Publishing.
- Malliga Subramanian, Rahul Ponnusamy, Sean Benhur, Kogilavani Shanmugavadivel, Adhithiya Ganesan, Deepti Ravi, Gowtham Krishnan Shanmugasundaram, Ruba Priyadharshini, and Bharathi Raja Chakravarthi. 2022. [Offensive language detection in Tamil Youtube comments by adapters and cross-domain knowledge transfer](#). *Computer Speech Language*, 76:101404.
- Praveenkumar Vijayakumar, Prathyush S, Aravind P, Angel S, Rajalakshmi Sivanaiah, Sakaya Milton Rajendram, and Mirnalinee T T. 2022. [SSN\\_ARMM@LT-EDI -ACL2022: Hope speech detection for equality, diversity, and inclusion using ALBERT model](#). In *Proceedings of the Second Workshop on Language Technology for Equality, Diversity and Inclusion*, pages 172–176, Dublin, Ireland. Association for Computational Linguistics.